

Reference-based publication networks with episodic memories

ANTHONY F. J. VAN RAAN

Centre for Science and Technology Studies, Leiden University, Leiden (The Netherlands)

In this paper we report first results of our study on network characteristics of a reference-based, bibliographically coupled (BC) publication network structure. We find that this network of clustered publications shows different statistical properties depending on the age of the references used for building the network. A remarkable finding is that only the network based on all references within publications is characterized by a degree distribution with a power-law dependence. This structure, which is typical for scale-free networks, disappears when selecting references of a specific age for the clustering process. Changing the publication network as a function of reference age, allows ‘tuning through the episodic memory’ of the nodes of the network. We find that the older the references, the more the network tends to change its structure towards a more exponential degree distribution.

Introduction

References are important characteristics of a publication. We studied linkages and clustering of publications from the year 2001 with help of ‘bibliographic coupling’ (BC) on the basis of their references (i.e., citations given to earlier publications) and measured the characteristics of the emerging network structure.

In bibliographic coupling, two articles are linked if they have at least one reference in common. Thus, a larger part of the scientific literature is structured by a network of interlinked publications that are often grouped in clusters. The BC structure is a rather unorthodox type of citation-based publication network in which *recent* literature – in this case publications of 2001 – is structured in terms of clusters on the basis of co-referencing, whereas in co-citation analysis (CC) *older* literature, namely references in 2001-publications are structured in terms of clusters on the basis of their co-citing papers. In other words, in BC the ‘nowadays landscape’ of scientific literature is created on the basis of their memories to older literature, and in CC a landscape of ‘older literature’ is created reflected as it were from nowadays publications.

Received November 30, 2004

Address for correspondence:

ANTHONY F. J. VAN RAAN

Centre for Science and Technology Studies, Leiden University

Wassenaarseweg 52, P. O. Box 9555, 2300 RB Leiden, The Netherlands

E-mail: vanraan@cwts.leidenuniv.nl

0138–9130/US \$ 20.00

Copyright © 2005 Akadémiai Kiadó, Budapest

All rights reserved

In the usual citation networks studied so far^{1–4} the nodes (or: vertices) are published articles and a directed link (or: edge) from article A to a previously published article B indicates that A cites B, i.e., article A gives a reference to article B. Measurement of the number of times a publication (node of the network) is cited, yields the ‘incoming’ degree distribution.⁵ Thus, the degree distribution $P(k)$ gives the probability that a randomly selected node has k links. The degree distribution is a kind of stationary (a time-independent) measure of the network.⁶

In a randomly wired network (random graph model) the nodes have a uniformly distributed probability to connect. The probability that a node has k links, and with that the degree distribution of these random networks, then follows a Poisson-distribution. An important characteristic of real networks, however, is local clustering which means that network-structures are more complex than simple randomly wired networks. It appears that many real networks are scale-free, i.e., their degree distribution follows a power law. Large networks may self-organize into such a scale-free state. Scale-free means that a functional form $f(x)$ remains unchanged under rescaling of a variable x , which means $f(ax) = bf(x)$. The solutions to this general equation are always power law forms.

The functionality of the network heavily depends on the type of distribution. Hence, characteristics such as degree distributions are not just of statistical interest. These distributions describe the structure of networks, and structure is directly related to important features of a network such as signal-propagation speed.⁷ An extensive overview of the statistical mechanics of complex networks is given by Albert and Barabási.⁸

Our study aims at finding statistical properties, such as degree distribution, path lengths, connectivity distributions and dynamical aspects that characterize the structure and behaviour of BC publication networks. In this paper we focus on phenomena related to the degree distributions. We also aim to understanding the meaning of these properties.

We take an analogy with scientific collaboration networks. In scientific collaboration networks, two nodes (authors) are linked if they coauthored one or more publications. These ‘co-author’ networks recently represent a kind of archetype example of a complex evolving network (see Refs 9–11 on static properties and Ref. 6 on dynamical properties). Therefore, we illustrate the analogy of our BC publication network with the scientific collaboration network in the following scheme:

Authors have publications	Publications have references
Publications may have more than one author	References may appear in more than one publication
These authors are called co-authors	These publications are called bibliographically coupled (BC) publications
An author may have s co-authors	A publication may have s BC publications
We call this a collaboration cluster of size s	We call this a BC publication cluster, size s
In a scientific collaboration network, authors are the nodes and co-authorship establish the links	In a scientific BC publications network, publications are the nodes and bibliographic coupling establish the links
Number of publications per author	Number of references per publication
Number of authors per publication	Number of publications per reference (this means in fact the number of citing papers in a year to a specific reference)

There is an interesting difference between the two above networks. In the scientific collaboration network, or ‘co-author’ network, a publication functions as an ‘affinity characteristic’ of authors: it is the element which causes scientists to cluster as co-authors. Publications (the ‘clustering elements’) and authors (the ‘clustered elements’) are, however, completely different entities in nature. In the BC publication network, a reference functions as an ‘affinity characteristic’ of publications: it is the element which causes publications to cluster as BC co-publications. But references are publications themselves, so in our network the ‘clustering elements’ and the ‘clustered elements’ are the same things in nature. It is like clustering people on the basis of their parents, grandparents, and further forefathers.

In the scientific collaboration network the co-authors know each other, they form a social structure of personal relations. In the BC publication network, the co-publications ‘do not know each other’, they just share one or more references, and in that sense this network looks like a large consumer system in which we find clusters of consumers (who do not necessarily know each other) sharing an interest in a specific groups of products or of services. In this study, the ‘shared interest’ is a research theme or research area.

Basic principles of BC publication networks

Our network consists of linked 2001-publications. These publications are connected by their referencing characteristics. This type of clustering is called *bibliographic coupling* (BC) as opposed to *co-citation coupling* (CC). The history of co-citation study

bibliographic coupling studies goes back to Kessler in the early 1960's. An extensive overview of the work of Kessler is given by De Solla Price in his pioneering work 'Networks of Scientific Papers'.¹²

We first explain the main lines of our method. We define a publication as a function of its references. As a simple example we take a small publication data set in which we have four publications p_1, p_2, p_3 and p_4 , and five references r_1, r_2, r_3, r_4 , and r_5 . In bibliometric language, the p_i are the *citing* papers, and the r_i the *cited* papers. Say p_1 contains all references r_1 to r_5 ; p_2 contains r_1, r_3 and r_4 ; p_3 has only r_1 and r_4 as a reference, and p_4 has none of the five references in its reference list. We can now construct a publication-to-reference matrix P :

	r_1	r_2	r_3	r_4	r_5
p_1	1	1	1	1	1
p_2	1	0	1	1	0
p_3	1	0	0	1	0
p_4	0	0	0	0	0

We observe that p_1 is bibliographically coupled to p_2 via r_1 (and also via r_3 and r_4 , but one 'link' is sufficient to have a bibliographic coupling). It is clear however that this number of links – three in this case – can be used as a measure of strength of the bibliographic coupling between two articles. Also, p_1 is coupled to p_3 (via r_1 , or via r_4), but not to p_4 . Thus, p_1, p_2 and p_3 form a BC-cluster, in which p_1 has two BC 'co-publications' (and the same is true for p_2 and p_3). Notice that p_1, p_2, p_3 and p_4 have together 10 references, which, however, represent 5 different cited articles.

Pre-multiplication of matrix P with its transpose P^T yields the (symmetric) reference-correlation matrix:

$$C = P^T \times P \tag{1}$$

which is in our example:

3	1	2	3	1
1	1	1	1	1
2	1	2	2	1
3	1	2	3	1
1	1	1	1	1

The *diagonal* values (printed in bold face) of this matrix C , $c(i,i) = c(i)$, indicate the number of times a specific cited publication (reference) is mentioned in the total set of publications. This is 3 times for r_1 , 1 time for r_2 , and so on. Thus, the matrix diagonal represents the *occurrences* of each reference, or, in other words, the number of citations given by the set of publications to each cited publication. This means that the distribution function of $c(i)$ gives the *in-degree* distribution of the publication set: the

number of (citing) publications per cited publication (reference). We discuss this in-degree distribution $N(c)$ for our empirical data in the next section.

The *off-diagonal* values give the *co-occurrences*, for instance r_1 and r_4 (value printed in italics) are mentioned 3 times together in publications of the set, namely in p_1 , p_2 , and p_3 , but for this information we have to go back to the original matrix \mathbf{P} as in matrix \mathbf{C} all ‘direct’ information on the *publications* is ‘lost’.

These co-occurrences of references are the basis of what is called in bibliometric studies ‘*co-citation* linkage clustering’ (CC). For instance: r_1 and r_4 are linked with ‘strength’ 3, but r_4 is also linked to r_3 with strength 2, and so on.

If we now take the mirrored matrix multiplication of Eq. 1, i.e., post-multiplication of our original matrix with its transpose, we get the (symmetric) publication-correlation matrix:

$$\mathbf{R} = \mathbf{P}(\mathbf{r}) \times \mathbf{P}^T(\mathbf{r}) \quad (2)$$

which is in our example:

$$\begin{array}{cccc} \mathbf{5} & 3 & 2 & 0 \\ 3 & \mathbf{3} & 2 & 0 \\ 2 & 2 & \mathbf{2} & 0 \\ 0 & 0 & 0 & \mathbf{0} \end{array}$$

The *diagonal* values (again printed in bold face) of this matrix \mathbf{R} , $r(i,i) = r(i)$, indicate the number of references in each publication. This is 5 for p_1 , 3 for p_2 , and so on. This means that the distribution function of $r(i)$ gives the *out-degree* distribution of the publication set: the number of references (cited publications) per (citing) publication. We discuss this out-degree distribution $N(r)$ for our empirical data also in the next section.

The *off-diagonal* values give the number of references shared by any two publications, for instance p_1 and p_3 (value printed in italics) share 2 references (namely in r_1 and r_4 , but also for this information we have to go back to the original matrix \mathbf{P} as in matrix \mathbf{R} all ‘direct’ information on the *references* is ‘lost’). Thus, \mathbf{R} provides the *strengths* of the links between each possible publication pair within the set and can be used for calculating ‘distances’ between publications.

The structure of the BC-coupled publications can be represented in a bipartite graph, and the resulting network as a simple projection of that graph, as given in Figure 1, together with the strengths of the links as follow from matrix \mathbf{R} . We observe that $\sum_{i,j(i \neq j)} r_{ij}$ gives the total number of links to node (publication) p_i , for instance for p_1 this number is 5, for p_2 it is also 5, and for p_3 it is 4. The above discussion shows that co-citation linking and bibliographic coupling are mathematically related by simple matrix algebra.

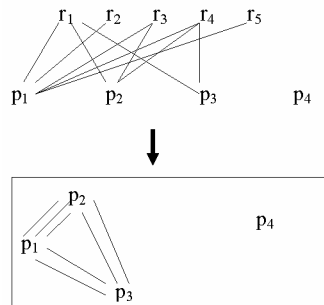


Figure 1. Graph of the BC network example

Results and discussion

Characteristics of the BC-network

We analyzed all 2001 publications in the complete set of citation indexes*, which totals to 1,099,017. The following characteristics of the entire data set were studied.

Number of references per publication. The number of references per publication is an out-degree measure (analogous to the number of publications per author in the case of scientific collaboration), i.e., the distribution $N(r)$ of 2001-publications as a function of the number of references, see Figure 2. For instance, we find that there are 43,364 publications in 2001 having just 1 reference, and one publication having 2,301 references.

This distribution function $N(r)$ appears to be a curve consisting of two power-law regimes. The tail, i.e., the higher- r part (from about $r = 40$) of the distribution, follows a steep power-law. This is the part of the distribution where a considerable amount of the publications are review articles. We find a power-law decay with exponent approximately -3.7 . The low- r part (particularly for r between 1 and 10) is quite flat and follows a very slowly decreasing power-law. Here most of the publications will be shorter articles such as letters, notes. Apparently there is not much difference in the (relatively small) number of references for these types of publications. Taking into account only those references that are themselves articles covered by the ISI databases and published after 1980 (in order to get a pair-wise match of a reference with a source article in our database), we find a power-law decay of approximately -3.5 .

* The Science Citation Index (SCI), the Social Science Citation Index (SSCI), the Arts & Humanities Citation Index (AHCI), and all 'specialty' indexes such as Neurosciences, Biochemistry and Biotechnology, etc., published by the Institute of Scientific Information (ISI), now Thomson Scientific, in Philadelphia.

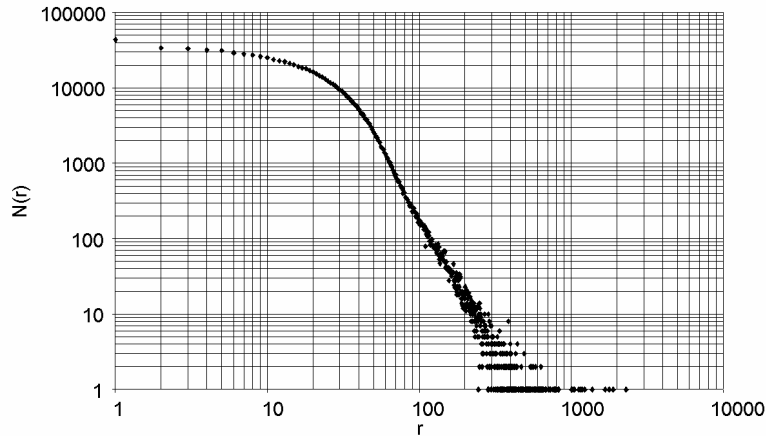


Figure 2. Number of references per publication, 2001

In his seminal paper, De Solla Price¹² reports a power-law behaviour in the tail of the distribution with an exponent approximately -2.0 for the out-degree of the citation network ('incidence of references'). He used the references of papers published in 1961. Also the flat low- r part was observed. It is not yet clear why we find a much steeper decay than in the work of De Solla Price. Possibly reference characteristics of publications have been changed in the last 40 years since Price's observations. Therefore, we are currently investigating the number of references per publication as a function of publication year with the data available in our bibliometric data-system, ranging from 1980 up till now.

Using the analogy of our BC-network with the scientific collaboration network, we may compare the number of references per publication with the number of publications per author. For this latter distribution, Newman⁹ finds an exponentially truncated power law. He attributes this truncation to the finite time window of five years used in his study, which prevents authors from having a very large number of papers. For instance, the co-author degree distribution for the Los Alamos Archive can be described with a truncated power law $P(k) = C k^{-\tau} \exp(-k/\kappa)$, τ and κ are specific parameters, see Ref. 9. In the original work of Lotka,¹³ which shows a 'complete' power-law, a more 'life time' approach was taken and therefore such a window was not used. Alternative explanations for the truncation are based on specific growth models of networks^{4,6} or specific collaboration models.¹⁵

Number of citations to the references. Number of citations in 2001 (i.e., number of citing publications in 2001) to references, $N(c)$, is an *in-degree* measure, analogous to

the number of authors per paper in the case of scientific collaboration. Thus, we take all references in the 2001-publications, and analyze the number of times these references are cited by the 2001-publications. The results are shown in Figure 3. In 2001 15,301,841 references are given and these references represent in fact 4,876,752 'unique' references (i.e., cited articles). We find 2,230,575 *cited articles* that are cited only once in 2001-publications; 203,407 articles cited five times in 2001; 401 articles cited fifty times in 2001, and 1 article that is cited 3,484 times. This is, by definition, the most cited paper (from 1980 on) in 2001.

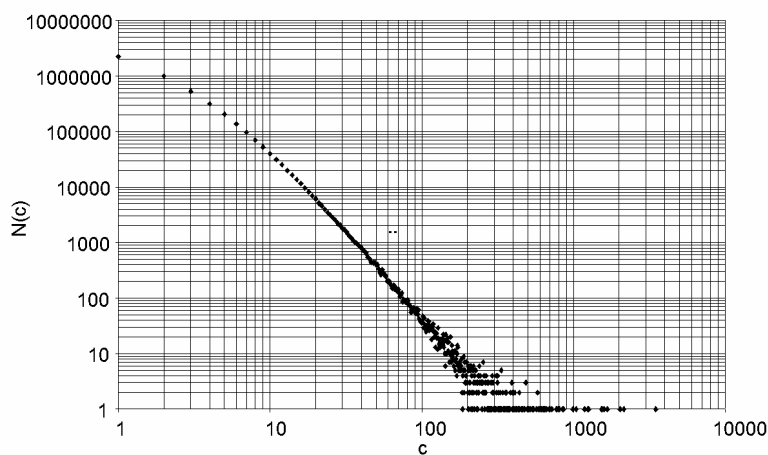


Figure 3. Number of citations in 2001 to references of 2001-publications

The citation distribution function shows for the major part a power-law decrease with a cut-off for lower c -values (c from 1 to about 10). We find a power-law decay with exponent approximately -3.1 .

This power-law distribution for citations is well-known. In the above mentioned pioneering work, De Solla Price also studied the in-degree distribution ('incidence of citations') and reports a power-law distribution with an exponent between about -2.5 and -3.0 .^{12,16} Also Naranan,¹⁷ using a subset of the data of Price, finds a power law exponent close to -3.0 .

On the basis of his observations, De Solla Price developed the model of 'cumulative advantage', building on Simon's work on the Matthew effect, i.e., the rich get richer (see for instance Bornholdt and Ebel¹⁸ and the original work of Simon¹⁹). In network-language, this phenomenon is a striking example of 'preferential attachment' (i.e., the probability for a node to obtain a new link increases with the number of links this node

already has) as in citation networks a new publication is likely to cite a well-known and thus mostly much-cited publication more than a less cited publication.^{5,20}

Measuring citation distributions are not so straightforward as often thought. There are quite different modalities of measurement. Redner¹ analysed citations to 1981-publications received in the years 1981–1997, thus the in-degree of publications of one year (1981). This is the case of a *fixed publication-year* (1981) followed by a wide ‘window’ of citation years. In our study we have a *fixed citation-year* (2001) and a wide window of preceding publications years. Also De Solla Price used a fixed (1961) citation year. Redner finds in his study that the asymptotic tail of the citation distribution appears to be described by a power law with exponent approximately -3 . But given the quite different behaviour of his distribution function for high versus low citation numbers, he suggests two different ‘citation regimes’, in the sense that there might be different underlying mechanisms and thus different statistical features between less-cited (exponential behaviour) and highly-cited papers (power law). We think, however, that further studies are necessary to investigate the effects caused by the difference in measuring modalities, as in our case the deviation from a power law behaviour for low citation number is less stronger than in the study of Redner.

Laherrère and Sornette²¹ and also Tsallis and De Albuquerque²² state that natural phenomena often exhibit a power-law followed by a significant curvature. They question whether these observed deviations form a power-law behaviour just simply result from finite-size effects or the existence of two regimes that are different in nature. They discuss models in which the distribution of citations of scientific papers can be fitted over the entire range of citation numbers with one single curve, so called stretched exponentials with the general form $f(x) \sim \exp\{- [(x/x_0)^\beta]\}$ (for the specific parameters x_0 and β see Ref. 22). This stretched exponential distribution is related to the Weibull distribution, see Ref. 32. Laherrère and Sornette²¹ apply these stretched exponentials (yielding ‘parabolic fractals’) to citations of highly cited physicists. Other examples are the size-distribution of cities. This is important, as in many cases claimed power laws clearly show a ‘parabolic’ effect for high k -values instead of a ‘real straight line’ (in a log-log plot), see for instance Figure 2b and also Figures 14 and 15 in Ref. 6. In these latter figures we see that this ‘parabolic’ effect is stronger with less nodes.

In a recent paper²³ we studied the in-degree distribution of about 15,000 chemistry publications published in The Netherlands in the period 1985–1993. Citations were counted in a modality again different from the two earlier mentioned. It is the modality often used in bibliometric analysis for evaluation purposes. For each publication year within the range 1985–1993, a 3-year window to receive citations after the year of publication is used. For instance, for publication year 1985 the citation window is 1986–1988, and so on. This measurement modality has the advantage of giving each publication year the same time period for receiving citations.

The resulting distribution function shows, for the larger number of citations, approximately a power law with an exponent of about -2.6 . We observe similar ‘parabolic’ deviations from the ‘ideal’ power law as discussed above: an inclination to ‘saturate’ for the lower citation values, as well as a cut-off for higher citation values. In contrast to fitting procedures as in the work of Laherrère and Sornette²¹ and of Tsallis and De Albuquerque,²² we developed a novel, *ab initio* theoretical model for the acquisition of citations by publications on the basis of a two-step competition process. Surprisingly, the result of this model is not the prediction of a power law behaviour for the citation distribution. We find a second order Bessel function. And even more surprisingly, this second order Bessel function describes the empirically measured distribution function very well, for the entire range of citation values. This would mean, that the mechanism of citation distribution only ‘mimics’ a scale-free (power law) behaviour. We are currently investigating the ability of our model to describe the in-degree distribution of the citation data in this study and refer for details and comparison with empirical results to Ref. 23.

Number of bibliographically coupled publications per publication. The number of bibliographically coupled publications, or ‘BC co-publications’, per publication, $N(s)$, is the main characteristic of our reference-based publication network system. Using again our analogy with the case of scientific collaboration, it is comparable with the number of coauthors (or: collaborators) per authors (see the scheme in the Introduction). A group of BC co-publications can be considered as a cluster of publications (just as a group of coauthors can be considered as a cluster), and therefore we also use the term ‘cluster size’ when dealing with the number of bibliographically coupled publications.

The distribution of the number of BC co-publications (BC-cluster size distribution) based on all references of the 2001-publications is presented in Figure 4a. In order to investigate an ‘episodic memory effect’ in the emergence of BC clusters, we also studied the distribution as a function of the age of the references, see Figures 4b, c, d, e, and f. This age-dependent measurement is performed by selecting from the total set of references five different subsets: the references (given in 2001-publications) with publication years 1998–2001 (the ‘youngest’ references, Figure 4b), publication years 1994–1997 (Figure 4c), publication years 1990–1993 (Figure 4d), publication years 1986–1989 (Figure 4e), and, finally, references with publication years 1985 and before (the ‘oldest’ references, Figure 4f), respectively.

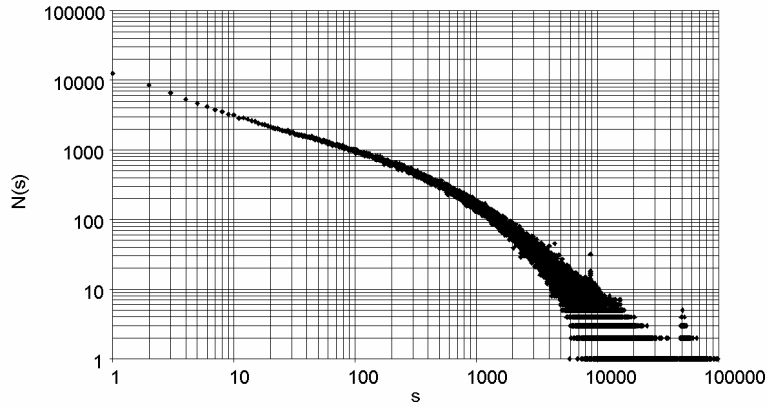


Figure 4a. Number of BC clusters, based on all refs.

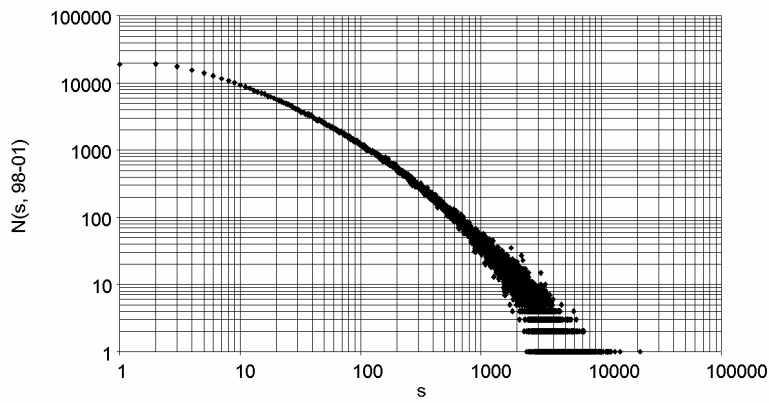


Figure 4b. Number of BC clusters, based on 1998–2001 refs.

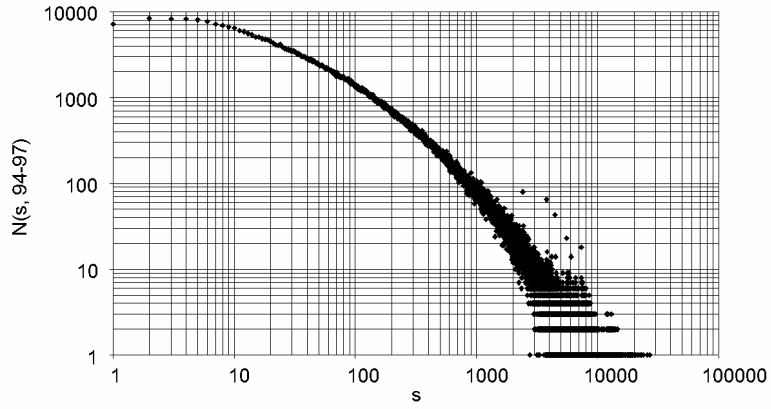


Figure 4c. Number of BC clusters, based on 1994–1997 refs.

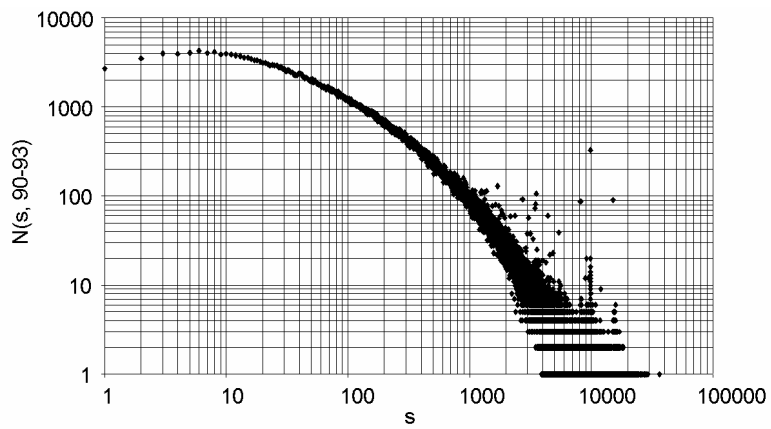


Figure 4d. Number of BC clusters, based on 1990–1993 refs.

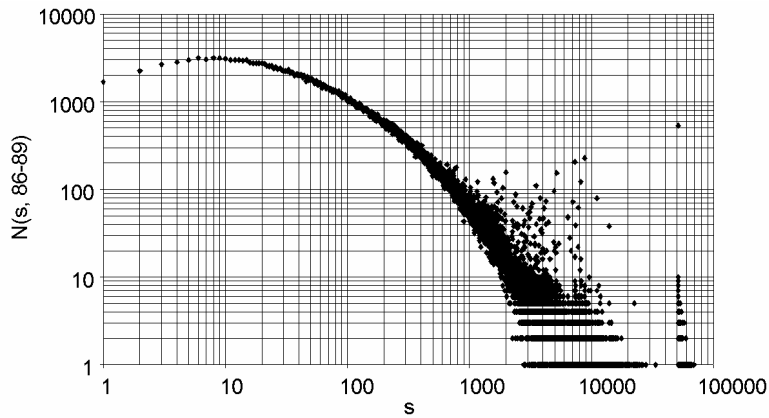


Figure 4e. Number of BC clusters, based on 1986–1989 refs.

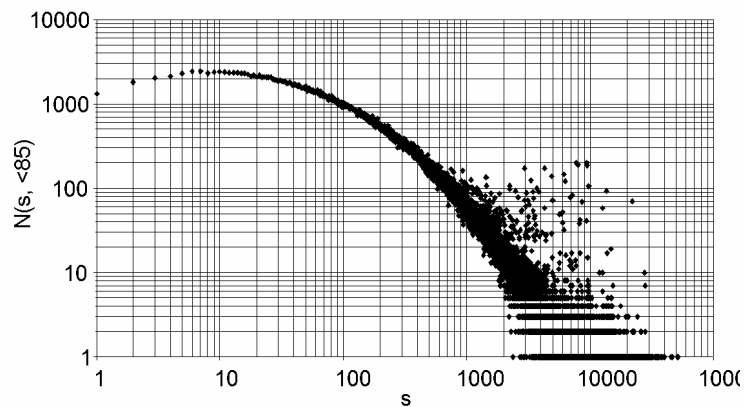


Figure 4f. Number of BC clusters, based on 1985 and earlier refs.

The analyses show remarkable results. In the case where publications are characterized by *all* their references (Figure 4a), we find a BC-cluster size distribution that is typical for a ‘scale-free’ network, i.e., with a power-law behaviour, with an exponential cut-off for cluster sizes above about 1,000. The *older* the references used to construct the BC network, the *stronger* the deviation of the distribution from a power-law toward a more exponential behaviour. This would mean that publications

characterized by just their oldest references, cluster in a much more random way than if the entire list of references is taken into account. In other words, clusters based on 'old memories' tend to be distributed more exponentially, which also means less small clusters, as can be observed in Figures 4b-f. As soon as the nodes 'rejuvenate', i.e., increase their 'short term memory', they tend to form a more power-law structured, i.e., scale-free network (much more small clusters). This tendency to deviate from a power law toward a more exponential behaviour in networks with 'aging of sites' is also observed in the model of Dorogovtsev and Mendes.²⁴

How can we explain this? The relatively old references are 'archival' and mostly much more *general* or 'classic' than the more recent references, which are typical field- or research theme-*specific*. Thus, these older references tend to link more parent publications, and the wiring of the BC network will therefore be more randomly distributed among the participating nodes, i.e., those 2001-publications having these relatively old references. Barabási et al.²⁵ show that in case of a growing network the degree distribution function has an exponential form in case of 'uniform attachment', i.e., the new node connects with equal probability to the nodes already present in the system, independent of the degree values of a node (no preferential attachment, their model B, see also Ref. 8). In our 'tuning' through the references, we more or less simulate a similar process in an otherwise static structure.

Smaller clusters are typical for research on very specific themes and in most cases these themes are very recent and thus characterized by relatively young references (the 'short-term memory of the system'). We indeed observe much more smaller clusters on the basis of 1998–2001 references (order of magnitude: 10,000) than in the cases of older references (order of magnitude: 1,000).

There is, however, an analytical problem. Since the *number* of references given by the citing publications is age-dependent, selection of increasingly older references also implies a choice for *increasingly less* references to characterize a publication. The largest group of references in publications concerns the most recent references, i.e., references to papers from 1998–2001. Also, it is obvious that the more references are included in the clustering process, the more larger clusters will be found. Indeed, we observe that the distribution based on the entire reference lists (Figure 4a) contains the largest clusters.

In order to distinguish between 'time-dependent' and 'less references', we again performed the BC wiring process, but now with removing randomly 10% of the references. As an example of this data manipulation for the BC clustering process with the oldest references (i.e., from 1985 an earlier) is given in Figure 5. We immediately observe that there is no significant difference in the shape of the distribution as compared with the 'complete' data shown in Figure 4f. Thus, we conclude that the topological changes reported in this paper are indeed due to time-dependent effects.

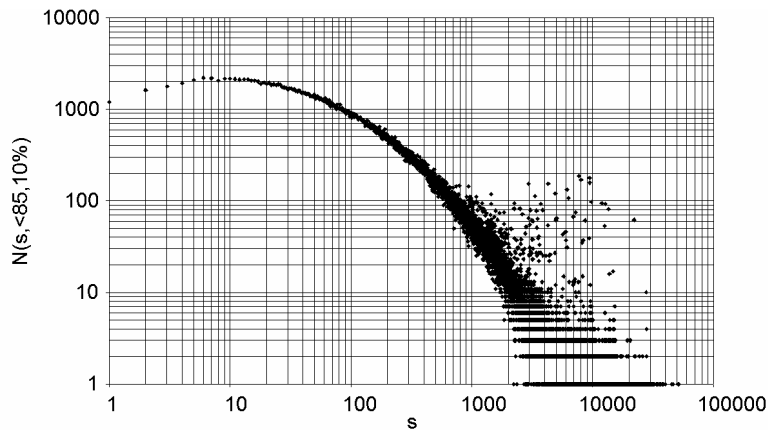


Figure 5. Number of BC clusters, based on 1985 and earlier refs, minus 10% random

It is fascinating to observe that *only* on the basis of the entire reference lists, publications create a BC network with – at least up to large cluster sizes – a scale-free, power-law behaviour. We conclude from this observation that the ‘affinity’ of publications with other publications is only optimal in the case of complete reference lists. If one deletes a part of a publication’s reference list, the publication is not anymore ‘what it is’, not fully characterized (remember that in the BC process a publication is as it were represented by its set of references). As specific ‘affinity’ leads to ‘preferential attachment’ (just as in molecular attachment processes), and as this preferential attachment is a strong condition for scale-free behaviour of networks, it is plausible to say that only in the case of BC networks based on the set complete of references we will find a power-law distribution function.

Removing references is a kind of imposing constraints to the publications, which are the nodes of the BC network. From earlier work we know that preferential attachment can be hindered by such constraints. Barabási and Albert²⁰ show the influence of constraints, demonstrating the difference between ‘physical’ networks (co-author networks, electrical power grids) which clearly ‘suffer’ from constraints, and ‘virtual’ networks such as the Worldwide Web where such physical constraints do not play a role (see also Refs 5, 8, 15). Amaral et al.²⁶ suggest that these constraints may determine the emergence of different classes of networks, and we believe that this effect is visible in our age-dependent ‘tuning’ through the references.

Concluding remarks

An important aspect of real-world networks is their growth.²⁷ Currently we are investigating the properties of the growth of our BC publication network. This network grows simply by adding the next year of publications, which would be 2002 in our case. Updating could also be done on a weekly or monthly base. Thus, for our network dynamical evolution is explicitly available. This growth process introduces an intriguing phenomenon. A smaller, but still considerable fraction of the references of these added 2002-publications are references to 2001-publications, which are nodes in the network.

But this does not mean that a new link is created between the new 2002-publication and the 2001-publication. A link between a new node (2002-publication) and an existing node (2001-publication) is only created if a reference of the new 2002-publication is the same as a reference in the 2001-publication, for instance both publications refer to an article published in 1999. In this rather curious way, even *clustering of old, unconnected nodes* is possible. For instance, 2001-publication p_1 has five references $r_1, r_2, r_3, r_4,$ and r_5 , while 2001-publication p_4 has none of these five references but contains other references, say, r_6 . In the 2001-network, p_1 and p_4 are *not linked* as they do not share any reference. If in 2002 a 'new' publication p_5 contains one or more of the p_1 references, for instance r_2 and r_3 , and also reference r_6 , this new publication p_5 will establish in the extended 2001 + 2002 network a link with the old nodes p_1 and p_2 . In other words, p_5 forms a BC cluster with p_1 and p_2 . Notice however that a *direct link* between the old nodes p_1 and p_2 is not created.

In general, for a newly added publication it is likely that one or several of its references will bibliographically couple this new publication to older publications that already have a large number of BC co-publications. Thus, in a growing BC network the 'older' nodes increase their connectivity leading to a reinforcement of preferential attachment. For a discussion of the measurement of preferential attachment in evolving networks, see Refs 5, 28.

Mossa et al.²⁹ make a connection with our earlier works^{30,31} on the growth of scientific literature by considering the situation in which new nodes are not processing information from a *constant fraction* of existing nodes, but from a *constant absolute number* of nodes. In other words, as the network grows, the new nodes are only able to process information from about a relatively small fraction of existing nodes. This model is plausible for networks that have grown to a very large size, for instance the scientific literature. They conclude that the above process reinforces preferential attachment, and with that, clustering which is in fact similar to fragmentation. Therefore, in our current work we investigate whether this is indeed the case for the BC network, which means, more concretely, whether the distribution function of number of references per paper does not change significantly as a function of time. Another interesting approach is to

replace the increasingly more *distant* reference windows (as in Figures 4a-f) by an increasing *wider* reference window, as to see how the shapes of the distributions will change. Instead of size-distributions, ranking distributions can be used to investigate more precisely the power-law character of the tail of the distributions. Furthermore, we will study in more detail how our citation distribution model discussed above could help to construct a theoretical framework to better describe the behaviour of processes taking place on networks.

*

The author thanks Peter Negenborn for his extensive data-analytical and programming work and the referees for several stimulating observations and suggestions.

References

1. S. REDNER, How popular is your paper? An empirical study of the citation distribution. *Eur. Phys. J. B4*, (1998) 131–134.
2. A. VAZQUEZ, Statistics of citation networks. *E-print arXiv: cond-mat/0105031*, (2001).
3. K. KLEMM, V. M. EGUÍLUZ, Highly clustered scale-free networks. *Physical Review E*, 65 (2002a) 036123.
4. M. E. J. NEWMAN, The structure and function of complex networks. *E-print arXiv: cond-mat/0303516*, (2003).
5. S. N. DOROGOVTSSEV, J. F. F. MENDES, *Advances in Physics*, 51 (2002) 1079–1187.
6. A.-L. BARABÁSI, H. JEONG, Z. NÉDA, E. RAVASZ, A. SCHUBERT, T. VICSEK, Evolution of the social network of scientific collaborations. *Physica A*, 311 (2002) 590–614.
7. D. J. WATTS, S. H. STROGATZ, Collective dynamics of ‘small-world’ networks. *Nature*, 393 (1998) 440–442.
8. R. ALBERT, A.-L. BARABÁSI, Statistical mechanics of complex networks. *Rev. Mod. Phys.*, 74 (2002) 47–97.
9. M. E. J. NEWMAN, Scientific collaboration networks. I. Network construction and fundamental results. *Physical Review E*, 64 (2001a) 016131.
10. M. E. J. NEWMAN, Scientific collaboration networks. II. Shortest paths, weighted networks, and centrality. *Physical Review E*, 64 (2001b) 016132.
11. M. E. J. NEWMAN, The structure of scientific collaboration networks. *Proc. Nat. Academy of Sciences*, 98 (2001c) 404–409.
12. D. J. DE S. PRICE, Networks of scientific papers. *Science*, 149 (1965) 510–515. See his ref. 3 for the papers of M. M. Kessler on bibliographic coupling.
13. A. J. LOTKA, The frequency distribution of scientific productivity. *J. Washington Acad. Sci.*, 16 (1926) 317–323.
14. P. L. KRAPIVSKY, S. REDNER, F. LEYVRAZ, Connectivity of growing random networks. *Phys. Rev. Lett.*, 85 (2000) 4629–4632.
15. R. ALBERT, A.-L. BARABÁSI, Topology of evolving networks: local events and universality. *Phys. Rev. Lett.*, 85 (2000) 5234–5237.
16. D. J. DE S. PRICE, *J. Amer. Soc. Inform. Sci. (JASIS)* 27 (1976) 292–306.
17. S. NARANAN, Power law relations in science bibliography- a self-consistent interpretation. *J. of Documentation*, 27 (1971) 83–97.
18. S. BORNHOLDT, H. EBEL, World Wide Web scaling exponent from Simon’s 1955 Model. *Phys. Rev. E*, 64 (2001) 035104.

19. H. A. SIMON, On a class of skew distribution functions. *Biometrika*, 42 (1955) 425–440.
20. A.-L. BARABÁSI, R. ALBERT, Emergence of scaling in random networks. *Science*, 286 (1999) 509–512.
21. J. LAHERRÈRE, D. SORNETTE Stretched exponential distributions in nature and economy: “fat tails” with characteristic scales. *Eur. Phys. J. B*, 2 (1998) 525–539.
22. C. TSALLIS, M. P. DE ALBUQUERQUE, Are citations of scientific papers a case of nonextensivity? *Eur. Phys. J. B*, 13 (2000) 777–780.
23. A. F. J. VAN RAAN, Two-step competition process leads to quasi power-law income distributions. Application to scientific publication and citation distributions. *Physica A*, 298 (2001) 530–536.
24. S. N. DOROGOVTSSEV, J. F. F. MENDEZ, Evolution of networks with aging of sites. *Phys. Rev. E*, 62 (2000) 1842–1845.
25. A.-L. BARABÁSI, R. ALBERT, H. JEONG, Mean-field theory for scale-free random networks. *Physica A*, 272 (1999) 173–187.
26. L. A. N. AMARAL, A. SCALA, M. BARTHÉLÉMY, H. E. STANLEY, Classes of small-world networks. *Proc. Nat. Academy of Sciences*, 97 (2000) 11149–11152.
27. K. KLEMM, V. M. EGUÍLUZ, Growing scale-free networks with small-world behavior. *Physical Review E*, 65 (2002b) 057102.
28. H. JEONG, Z. NÉDA, A.-L. BARABÁSI, Measuring preferential attachment in evolving networks. *Europhys. Lett.*, 61 (2003) 567–572.
29. S. MOSSA, M. BARTHÉLÉMY, H. E. STANLEY, L. A. N. AMARAL, Truncation of power law behavior in ‘scale-free’ network models due to information filtering. *Phys. Rev. Lett.*, 88 (2002) 138701.
30. A. F. J. VAN RAAN, On growth, ageing and fractal differentiation of science. *Scientometrics*, 47 (2000) 347–362.
31. A. F. J. VAN RAAN, Fractal dimension of co-citations. *Nature*, 347 (1990) 626.
32. T. NABESHIMA, Y.-P. GUNJI, Zipf’s law in phonograms and Weibull distribution in ideograms: comparison of English with Japanese. *BioSystems*, 73 (2004) 131–139.